Coarse-To-Fine Incremental Few-Shot Learning - Appendix

Xiang Xiang¹, Yuwen Tan¹, Qian Wan¹, Jing Ma¹, Alan L. Yuille², and Gregory D. Hager²

¹ MoE Key Lab of Image Processing and Intelligent Control, School of AI and Automation, Huazhong Univ. of Science and Tech., China ² Department of Computer Science, Johns Hopkins University, USA xex@hust.edu.cn

1 Full Proof of Proposition 1

Proposition 1 (Normalizing or freezing weights improves stability; doing both improves the most). Given Θ_a , if we only normalize weights of a linear FC classifier, we obtain Θ_b ; if we only freeze them, we obtain Θ_c ; if we do both, we obtain Θ_d . Then, $\mathcal{D}_d < \mathcal{D}_b < \mathcal{D}_a$ and $\mathcal{D}_d < \mathcal{D}_c < \mathcal{D}_a$. **Proof.** (1) Stability Degree of model Θ_a .

It is assumed that the training for all sessions will reach the minimum loss. For the training sample m in the 0-th session, the probability that m belongs to superclass is one, i.e., $p_{t,c_{super}}^m = 1$ and $p_{t,i}^m = 0 (i \neq c_{super})$. According to $p_i^m = \frac{\exp(o_i^m)}{\sum_{j=1}^I \exp(o_j^m)}$, the following conditions are satisfied,

$$\tilde{\mathbf{o}}_{c_{super}}^{(t)} = a(a \in \mathbb{R}), \tilde{\mathbf{o}}_{i}^{(t)} (i \neq c_{super}) = -\infty.$$
(1)

After training of T-th session has reached the minimum loss, $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = b(b \in \mathbb{R}), \tilde{\mathbf{o}}_{i}^{(T)} (i \neq c_{sub}) = -\infty$, then,

$$\mathcal{D}_{a} = \sum_{i} \left(\frac{\tilde{\mathbf{o}}_{i}^{(T)} - \tilde{\mathbf{o}}_{i}^{(t)}}{\tilde{\mathbf{o}}_{i}^{(t)}}\right)^{2} = \left(\frac{-\infty - a}{a}\right)^{2} + \left(\frac{b - (-\infty)}{-\infty}\right)^{2} = \infty.$$
(2)

(2) Stability Degree of model Θ_b .

Under the same conditions above, the following conditions are satisfied according to $p_i^m = \frac{\exp(\cos \theta_i^m)}{\sum_{j=1}^I \exp(\cos \theta_j^m)}$,

$$\tilde{\mathbf{o}}_{c_{super}}^{(t)} = 1, \tilde{\mathbf{o}}_{i}^{(t)} (i \neq c_{super}) = -1.$$
(3)

After training of *T*-th session has reached minimum loss, $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = 1$, $\tilde{\mathbf{o}}_{i}^{(T)}$ $(i \neq c_{sub}) = -1$, then the following holds:

$$\mathcal{D}_b = \sum_i \left(\frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}}\right)^2 = \left(\frac{-1-1}{1}\right)^2 + \left(\frac{1-(-1)}{-1}\right)^2 = 8.$$
(4)

(3) Stability Degree of model Θ_c .

2X. Xiang et al.

Compared with Θ_a , model Θ_c freezes weights of neurons corresponding to previously-seen classes. After training of T-th session has reached its minimum loss, $\tilde{\mathbf{o}}_{c_{super}}^{(T)} = a, \tilde{\mathbf{o}}_{c_{sub}}^{(T)} = \infty^+, \tilde{\mathbf{o}}_i^{(T)} (i \neq c_{super} \lor i \neq c_{sub}) = -\infty$, where $\infty^+ > \infty$ in order to offset the influence of $\mathbf{\tilde{o}}_{c_{super}}^{(T)}$, then,

$$\mathcal{D}_{c} = \sum_{i} \left(\frac{\tilde{\mathbf{o}}_{i}^{(T)} - \tilde{\mathbf{o}}_{i}^{(t)}}{\tilde{\mathbf{o}}_{i}^{(t)}}\right)^{2} = \left(\frac{\infty^{+} - (-\infty)}{-\infty}\right)^{2},$$

$$9 > \mathcal{D}_{c} > 4.$$
(5)

(4) Stability Degree of model Θ_d .

Compared with Θ_b , model Θ_d freezes weights of neurons corresponding to previously-seen classes. After training of T-th session has reached its minimum loss, $\tilde{\mathbf{o}}_{c_{super}}^{(T)} = 1, \tilde{\mathbf{o}}_{c_{sub}}^{(T)} = 1, \tilde{\mathbf{o}}_{i}^{(T)} (i \neq c_{super} \lor i \neq c_{sub}) = -1$, then,

$$\mathcal{D}_d = \sum_i (\frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}})^2 = (\frac{1 - (-1)}{-1})^2 = 4.$$
(6)

Comparing the stability degree of different models, we have $\mathcal{D}_{max} = \mathcal{D}_a$, $\mathcal{D}_{min} = \mathcal{D}_d$ and Θ_d is the most stable.

Full Analysis of Conjecture 2 $\mathbf{2}$

Conjecture 2 (Sufficient & necessary condition of no impact of freezing embeddingweights). $p \lor q \Leftrightarrow \neg r$ where p: classifier-weights are normalized, q: classifierweights are frozen, r: freezing embedding-weights improves the performance. The 'only if' part: $\neg p \land \neg q \Rightarrow r$

Analysis. We have 4 propositions that are all true according to our observations: $\neg p \land \neg q \to r$

- 2 3 $\begin{array}{c} p \land q \to \neg r \\ p \land \neg q \to \neg r \end{array}$
- (4) $\neg p \land q \rightarrow \neg r.$

They share a similar composition pattern. We summarize them as Table 1.

$P \mid Q \mid P \land Q \mid R \mid P \land Q \to R$						
$\begin{array}{c c} \neg p & \neg q \\ p & q \\ p & \neg q \\ \neg p & q \end{array}$	$ \begin{vmatrix} \neg p \land \neg q & r \\ p \land q & \neg r \\ p \land q & \neg r \\ p \land \neg q & \neg r \\ \neg p \land q & \neg r \\ \neg p \land q & \neg r \\ \neg p \land q & \neg r \end{vmatrix} $					

Table 1: Compound propositions.

Let us make p, q, r an realization of general propositions P: classifier-weights are normalized,

Q: classifier-weights are frozen,

R: freezing embedding-weights improves the performance, respectively. We want to construct a common proposition for all the four cases all to be true. Namely, we need to solve for a comopsitive proposition $\mathcal{C}(P,Q) \to R$ that satisfies the truth table with (1), (2), (3), (4) ordered top-down.

$P Q \mathcal{C}(P)$	$,Q) R \mathcal{C}(\mathcal{I})$	$\mathcal{P},\mathcal{Q}) \to R$
0 0	1	1
1 1	0	1
1 0	0	1
0 1	0	1

Table 2: A truth table that is not completed.

Note that $A \to B$ is 0 *iff* A is 1 and B is 0. Therefore, we want $\mathcal{C}(P,Q)$'s truth value of the 2, 3, 4 line never to be 1. Given the value pairs of P and Q, the only way to make that happen is to let $\mathcal{C}(P,Q)$ be $\neg P \land \neg Q \to R$, which is a solution that satisfies all four cases, and thus is always true.

$P Q \neg P \land \neg Q R \neg P \land \neg Q \to R$							
0	0	1	1	1			
1	1	0	0	1			
1	0	0	0	1			
0	1	0	0	1			

Table 3: The truth table is realized.

Namely, we have $\neg P \land \neg Q \Rightarrow R$, which is exactly Conjecture 2, $\neg p \land \neg q \Rightarrow r$, with a change of notations.

The 'if' part: $p \lor q \Rightarrow \neg r$.

Analysis. Given propositions (2), (3), (4), we will combine them and derive a logically-equivalent premise.

 $\begin{array}{l} (p \land q) \lor (p \land \neg q) \lor (\neg p \land q) \\ \Leftrightarrow (p \lor (p \land \neg q) \lor (\neg p \land q))) \land (q \lor (p \land \neg q) \lor (\neg p \land q)) \\ \Leftrightarrow ((p \lor (p \land \neg q) \lor \neg p) \land (p \lor (p \land \neg q) \lor q)) \\ \land ((q \lor (p \land \neg q) \lor \neg p) \land (q \lor (p \land \neg q) \lor q)) \\ \Leftrightarrow ((p \lor p \lor \neg p) \lor (p \lor \neg q \lor \neg p) \land (p \lor \lor q) \land (p \lor \neg q \lor q)) \\ \land ((q \lor p \lor \neg p) \land (q \lor \neg q \lor \neg p) \land (q \lor p \lor q) \land (q \lor \neg q \lor q)) \\ \land ((q \lor p \lor \neg p) \land (q \lor \neg q \lor \neg p) \land (q \lor p \lor q) \land (q \lor \neg q \lor q)) \\ \Leftrightarrow (1 \land 1 \lor (p \lor p \lor q) \land 1) \land (1 \land 1 \land (q \lor p \lor q) \land 1) \\ \Leftrightarrow (p \lor p \lor q) \land (q \lor p \lor q) \\ \Leftrightarrow (p \lor q) \land (p \lor q) \Leftrightarrow p \lor q. \end{array}$

4 X. Xiang et al.

Similarly, we can derive $(2) \lor (3) \lor (4)$ as $(p \land q \to \neg r) \lor (p \land \neg q \to \neg r) \lor (\neg p \land q \to \neg r)$ $\Leftrightarrow (p \land q) \lor (p \land \neg q) \lor (\neg p \land q) \to \neg r.$ With the premise replaced, we have

 $(2) \lor (3) \lor (4) \Leftrightarrow p \lor q \to \neg r,$

Given (2), (3), (4) are all always true. it holds that $p \lor q \to \neg r$ is always true, Namely, we have $p \lor q \Rightarrow \neg r$.

3 Insights from Fine-tuning

3.1 Embeddings need to be contrastively learned

As shown in Fig. 1, straightforward training on coarse labels does not help much the subsequent FSL on fine labels (now_acc at ~ 25%), while contrastive learning self-supervised by the fine cues does help (now_acc at ~ 35%). Thus, coarsely-trained embedding can be generalizable.

Fig. 2-left shows that freezing embedding-weight outperforms not freezing them. It implies the embedding space without any update is generalizable, and that, *if classifier-weights are not frozen, freezing embedding-weights helps.*

3.2 Freezing weights helps, surely for classifiers

Fig. 2-right implies that, *if classifiers weights are frozen, then freezing embedding-weights does not help.* Comparing left with right of Fig. 2, we find that freezing classifier-weights (right) outperforms not doing so (left), either freezing embedding-weights (circle) or not (triangle).



Fig. 1: Ablation study of contrastive learning when fine-tuning ResNet12 w/o IL. Left: w/o; right: w/. (CIFAR-100)



Fig. 2: Ablation study of freezing embedding-weights for fine-tuning a contrastive model. Left: when not freezing classifier-weights. Right: when freezing them.

4 Rethinking C2FSCIL and Knowe with More Results

Given a tree-like product catalog at Amazon.com, there is a class hierarchy per tree, as shown in Fig. 3. Furthermore, Fig. 4 presents the basic idea of Knowe. To learn over time (*i.e.*, sequential learning), it is suggested in [4,5] that neural networks can be limited by catastrophic forgetting (CF) just like Perceptron is unable to solve X-OR. Knowledge forgetting, or called catastrophic forgetting/interference is about a learner's memory (*e.g.*, LSTM) and is a result of the stability–plasticity dilemma regarding how to design a model that is sensitive to, but not radically disrupted by, new input [4,5]. Often, maintaining plasticity results in forgetting old classes while maintaining stability prevents the model from learning new classes, which may be caused by a single set of shared weights. Our setting requires a balance between coarse and fine performance unexplored by existing works, as shown in Fig. 5.



Fig. 3: Two examples of Amazon item catalog. Best seen on computer.



Fig. 4: Basic idea of Knowe. In base session we train Θ on \mathbb{D} to get $\Theta^{(0)}$. Per incremental session, $\Theta^{(t)}$ is trained on *C*-way *K*-shot support set $\mathbb{S}^{(t)}$ based on $\Theta^{(t-1)}$, $t \geq 1$ and then tested on any class seen in either \mathbb{D} or $\mathbb{S}^{(1)}, ..., \mathbb{S}^{(t)}$.

Method	Class hierarchy	Few-shot Learning I	ncremental Learning
LwF [3]			\checkmark
CEC [6]		\checkmark	\checkmark
ANCOR [2]	\checkmark	\checkmark	
IIRC [1]	\checkmark		\checkmark
C2FSCIL (Ours)	\checkmark	\checkmark	\checkmark

Table 4: Comparison of settings with related works.

References

1. Abdelsalam, M., Faramarzi, M., Sodhani, S., Chandar, S.: IIRC: Incremental Implicitly-Refined Classification. In: CVPR (2021)



Fig. 5: The stability-plasticity trade-off. Top-right is FT w/o IL; bottom-left represents most IFSL methods; bottom-right is our approach. (CIFAR-100)



Fig. 6: Knowe reaches a balance on BREEDS.

 Bukchin, G., Schwartz, E., Saenko, K., Shahar, O., Feris, R., Giryes, R., Karlinsky, L.: Fine-grained angular contrastive learning with coarse labels. In: CVPR (2021)

Fig. 7: The visualization of all confusion matrices of Knowe tested on the BREEDS living17 dataset.

- 3. Gidaris, S., Komodakis, N.: Dynamic few-shot visual learning without forgetting. In: CVPR (2018)
- McCloskey, M., Cohen, N.: Catastrophic interference in connectionist networks: the sequential learning problem. The Psychology of Learning and Motivation 24, 109– 164 (1989)
- 5. M.French, R.: Catastrophic forgetting in connectionist networks. Trends in Cognitive Sciences **3** (1999)
- Zhang, C., Song, N., Lin, G., Zheng, Y., Pan, P., Xu, Y.: Few-shot incremental learning with continually evolved classifiers. In: CVPR (2021)